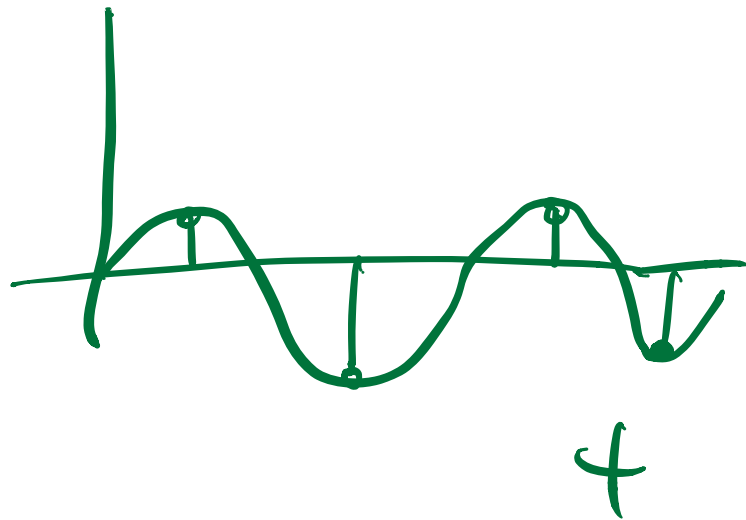
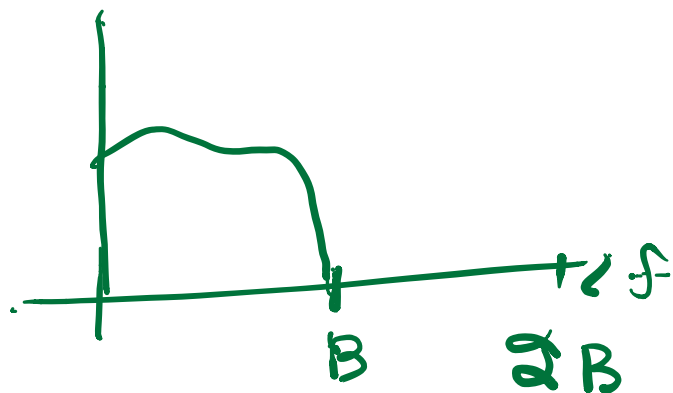


Sparsity Motivation



$$y = Fx$$

$\begin{bmatrix} \times \\ \times \\ \times \\ \times \end{bmatrix}$

time sig

$=$
 $\begin{bmatrix} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{bmatrix}$

DFT⁻¹

$\begin{bmatrix} \times \\ \times \\ \times \\ \times \end{bmatrix}$

Freq

$$[] = [] \begin{bmatrix} - \\ - \\ - \\ - \end{bmatrix}$$

↑
any matrix

underdetermined

x is sparse

1 OMP greedy

2.) LASSO

$$\min_x \underbrace{\|y - Ax\|_2^2}_{LS}$$

SP

$$\lambda \|x\|_0$$

↑ NP

!

$$+ \lambda \|x\|_1$$

Projects

- **3-4** person groups
- Deliverables: Poster, Report & main code (plus proposal, midterm slide)
- Topics: your own or chose from **suggested topics / kaggle**
- **Week 3 groups** due to TA Nima. Rearrangement might be needed.
- **May 2** proposal due. TAs and Peter can approve.
- Proposal: One page: Title, A large paragraph, data, weblinks, references.
- Something **physical and data oriented**.
- **May ~16** Midterm slides. Likely presented in 4 subgroups (3TA+Peter).
- **5pm 6 June** Jacobs Hall lobby, final poster session. **Snacks**
- Poster, Report & main code. Report due Saturday 16 June.

Final Report

Poster on **June 6** from each group is mandatory. Upload poster as well.

For the Final project (Due Saturday 16 June). Delivery Dropbox request <2GB (details to follow).:

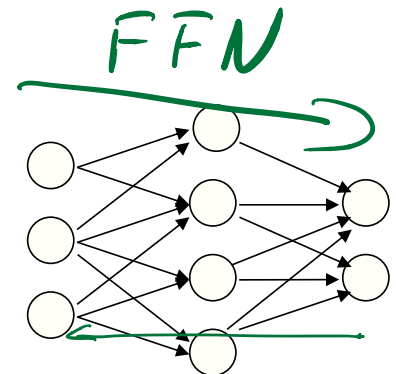
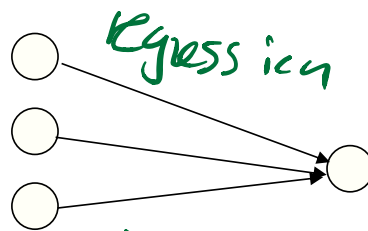
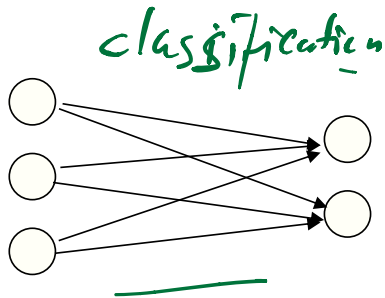
Deliver a code:

- Assume we have reasonable compilers installed (we use Mac OSX)
- Give instructions if any additional software should be installed.
- You can ask us to download a dataset. Or include it in this submission
- Don't include all developed codes. Just key elements.
- We should not have to reprogram your code.

Report

- The report should include all the following sections: Summary -> Introduction->Physical and Mathematical framework->Results.
- Summary is a combination of an abstract and conclusion.
- Plagiarism is not acceptable! When citing use “ “ for quotes and citations for relevant papers.
- Don't write anything you don't understand.
- Everyone in the group should understand everything that is written. If we do not understand a section during grading we should be able to ask any member of the group to clarify. You can delegate the writing, but not the understanding.
- Use citations. Any concepts which are not fully explained should have a citation with an explanation.
- Please be concise. Equations are good. Figures essential. Write as though your report is to be published in a scientific journal.
- Last year's reports are on class website. Especially good projects 2,11,12,13,

Lecture 8: Backpropagation



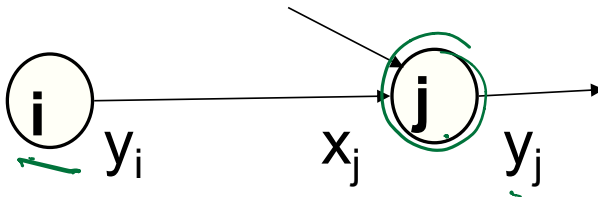
A difference in notation

- For networks with multiple hidden layers Bishop uses an explicit extra index to denote the layer.
- The lecture notes use a simpler notation in which the index denotes the layer implicitly.

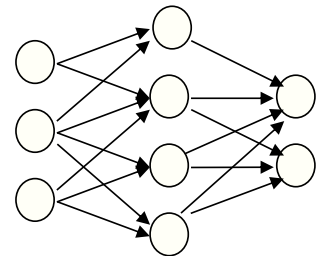
y is the output of a unit in any layer

x is the **summed** input to a unit in any layer

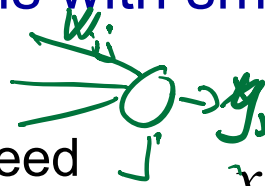
The index indicates which layer a unit is in.



Jeff Hinton



Non-linear neurons with smooth derivatives



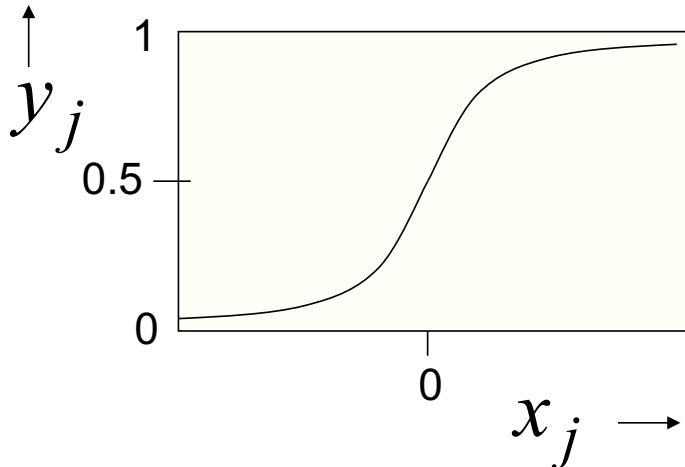
- For backpropagation, we need neurons that have well-behaved derivatives.

$$x_j = \underline{b_j} + \sum_i^I y_i w_{ij}$$

$$y_j = \frac{1}{1 + e^{-x_j}} \quad \checkmark$$

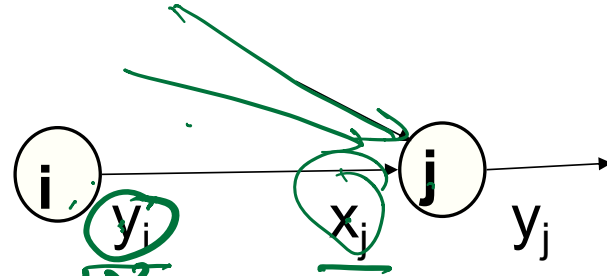
$$\frac{\partial x_j}{\partial w_{ij}} = y_i \quad \underline{\frac{\partial x_j}{\partial y_i} = w_{ij}}$$

$$\underline{\frac{dy_j}{dx_j} = y_j (1 - y_j)}$$



Backpropagation

- J nodes
- Observations t_j
- Predictions y_j
- Energy function $E = \frac{1}{2} \sum_j (y_j - t_j)^2$



- $\frac{dE}{dy_j} = y_j - t_j$

- $\frac{dE}{dx_j} = \frac{dE}{dy_j} \frac{dy_j}{dx_j} = (y_j - t_j) y_j (1 - y_j)$

- $\frac{dE}{dw_{ij}} = \frac{dE}{dx_j} \frac{dx_j}{dw_{ij}} = \frac{dE}{dx_j} y_i$

- $\frac{dE}{dy_i} = \sum_j \frac{dE}{dx_j} \frac{dx_j}{dy_i} = \sum_j \frac{dE}{dx_j} w_{ij}$

- $\frac{dE}{dx_i} = \frac{dE}{dy_i} \frac{dy_i}{dx_i} = \sum_j \frac{\partial E}{\partial x_j} w_{ij} \frac{dy_i}{dx_i}$

TensorFlow and Matlab training

- `tf.train.Optimizer`
- `tf.train.GradientDescentOptimizer`
- `tf.train.AdadeltaOptimizer`
- `tf.train.AdagradOptimizer`
- `tf.train.AdagradDAOptimizer`
- `tf.train.MomentumOptimizer`
- `tf.train.AdamOptimizer`
- `tf.train.FtrlOptimizer`
- `tf.train.ProximalGradientDescentOptimizer`
- `tf.train.ProximalAdagradOptimizer`

$$W_{new} = W_{old} - \eta \frac{\partial E}{\partial W}$$

Choose a Multilayer Neural Network Training Function

R2018a

It is very difficult to know which training algorithm will be the fastest for a given problem. It depends on many factors, including the complexity of the problem, the number of data points in the training set, the number of weights and biases in the network, the error goal, and whether the network is being used for pattern recognition (discriminant analysis) or function approximation (regression). This section compares the various training algorithms. Feedforward networks are trained on six different problems. Three of the problems fall in the pattern recognition category and the three others fall in the function approximation category. Two of the problems are simple "toy" problems, while the other four are "real world" problems. Networks with a variety of different architectures and complexities are used, and the networks are trained to a variety of different accuracy levels.

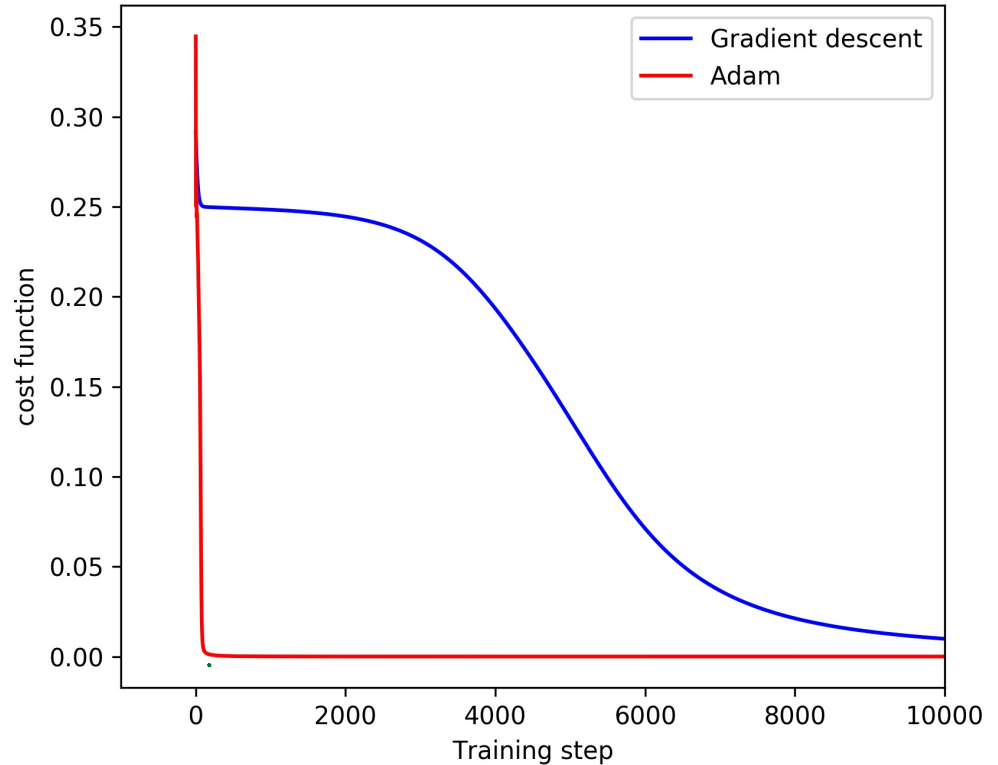
The following table lists the algorithms that are tested and the acronyms used to identify them.

Acronym	Algorithm	Description
LM	<code>trainlm</code>	Levenberg-Marquardt
BFG	<code>trainbfg</code>	BFGS Quasi-Newton
RP	<code>trainrp</code>	Resilient Backpropagation
SCG	<code>trainscg</code>	Scaled Conjugate Gradient
CGB	<code>traincgb</code>	Conjugate Gradient with Powell/Beale Restarts
CGF	<code>traincgf</code>	Fletcher-Powell Conjugate Gradient
CGP	<code>traincgp</code>	Polak-Ribière Conjugate Gradient
OSS	<code>trainoss</code>	One Step Secant
GDX	<code>traingdx</code>	Variable Learning Rate Backpropagation

The following table lists the six benchmark problems and some characteristics of the networks, training processes, and computers used.

Gradient results

Maybe from tensorflow with ocean acoustics data



Gradients

$$\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla E$$

- Batch

- $E(\mathbf{w}_k) = \sum_n^N E_n(\mathbf{w}_k)$

- Stochastic gradient descent:

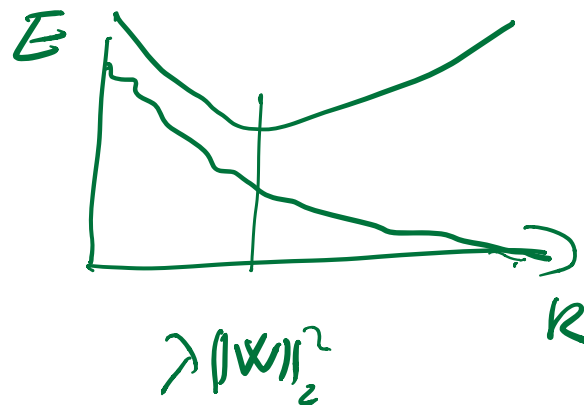
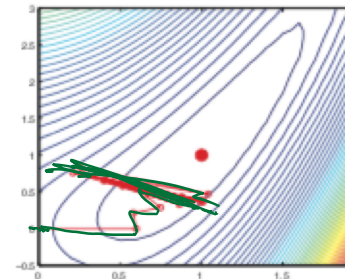
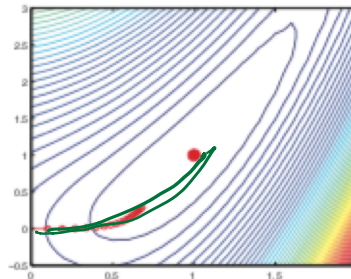
- Pick just one data sample n (of N)
 - $\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla E_n(\mathbf{w}_k)$
 - Less sensitive to global minima

- Momentum

- $\mathbf{w}_{k+1} = \mathbf{w}_k - \eta \nabla E + \underbrace{\mu(\mathbf{w}_k - \mathbf{w}_{k-1})}$

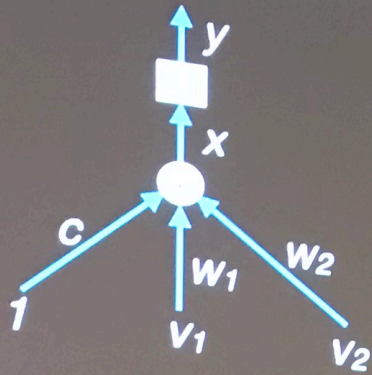
- Avoid Overfitting (regularization, dropout, early stopping)

- 100's PhD thesis on how to optimize. Always backpropagation.



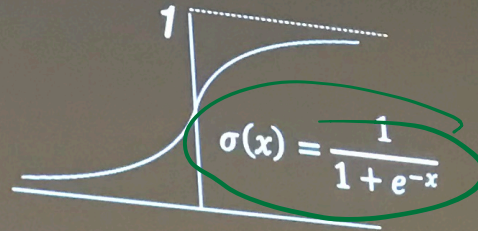
ICASSP 2018 Plenary

Stochastic Gradient Descent (SGD)



$$y = \sigma(x)$$

$$x = c + \mathbf{v}^T \mathbf{w}$$

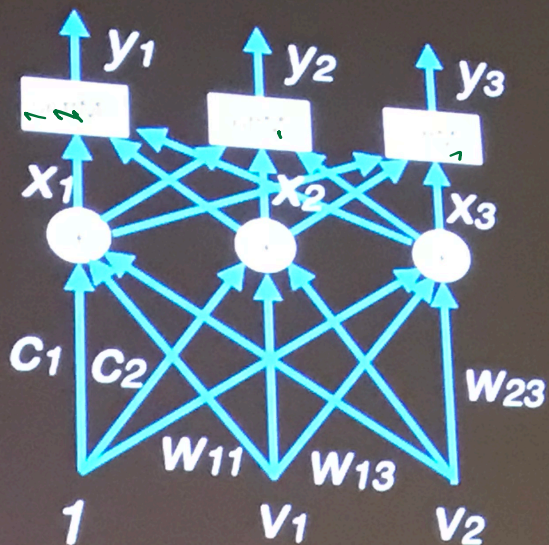


$$p(t|\mathbf{v}) = y^t(1 - y)^{1-t}$$

$$L = \ln p(t|\mathbf{v}) = t \ln y + (1 - t) \ln(1 - y)$$

$$\frac{\partial L}{\partial w_1} = \left(\frac{\partial L}{\partial y} \right) \left(\frac{\partial y}{\partial x} \right) \left(\frac{\partial x}{\partial w_1} \right) = (y - t)v_1$$

N-ary Classification



$$y_i = p(i|\mathbf{v}) =$$

$$\frac{e^{x_i}}{\sum_{l=1}^N e^{x_l}}$$

Softmax

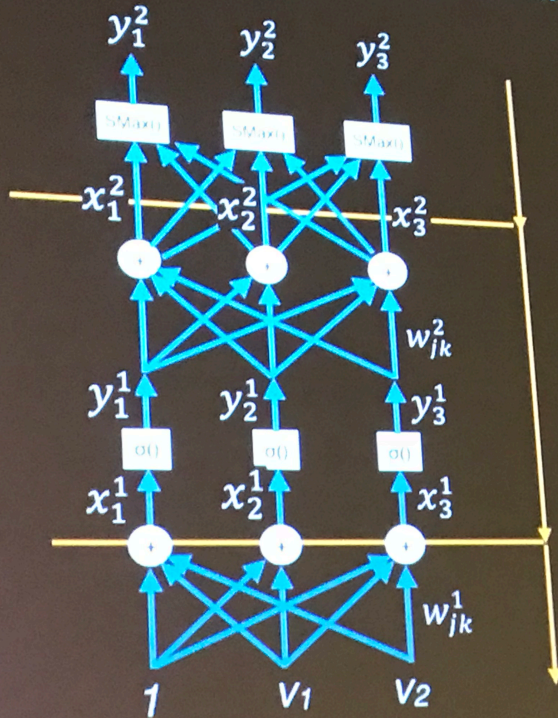
$$L = \sum_{i=1}^N t_i \ln y_i$$

$$w_{nj}^{(i)} = w_{nj}^{(i-1)} + \eta v_n^{(i-1)} (t_j - y_j^{(i-1)})$$

ICASSP 2018 Plenary

Perceptron Learning

Werbos, 1974; Rumelhart, Hinton, Williams 1986



Two-layers
2 input features
3 output labels

$$\nabla_n^2(m) = y_n^2(m) - t_n(m)$$

$$[w_{jn}^2]^{(i)} = [w_{jn}^2]^{(i-1)} - \eta \frac{1}{M} \sum_{m=1}^M y_n^1(m) \nabla_n^2(m)$$

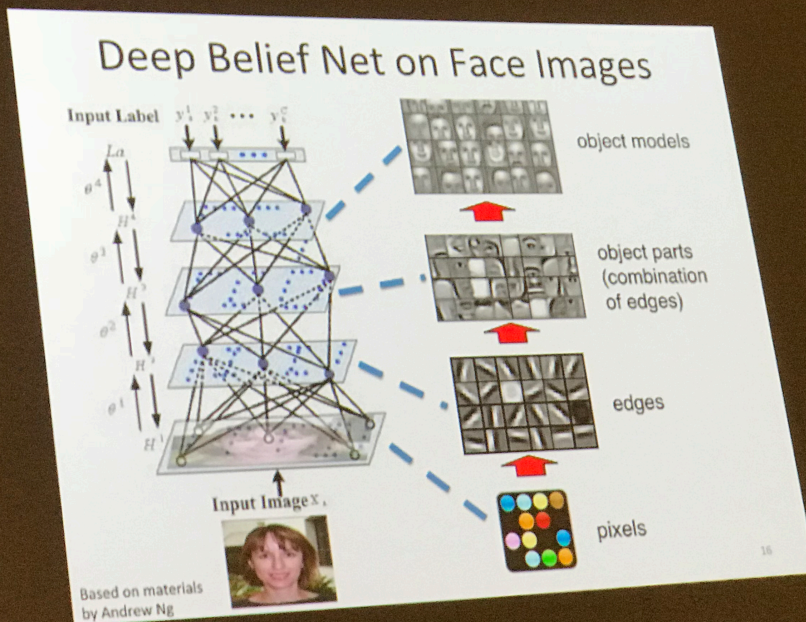
$$\nabla_n^1(m) = y_n^1(m)(1 - y_n^1(m)) \sum_{k=1}^N w_{nk}^2 \nabla_k^2(m)$$

$$[w_{jn}^1]^{(i)} = [w_{jn}^1]^{(i-1)} - \eta \frac{1}{M} \sum_{m=1}^M v_j(m) \nabla_n^1(m)$$

backprop
Min

ICASSP 2018 Plenary

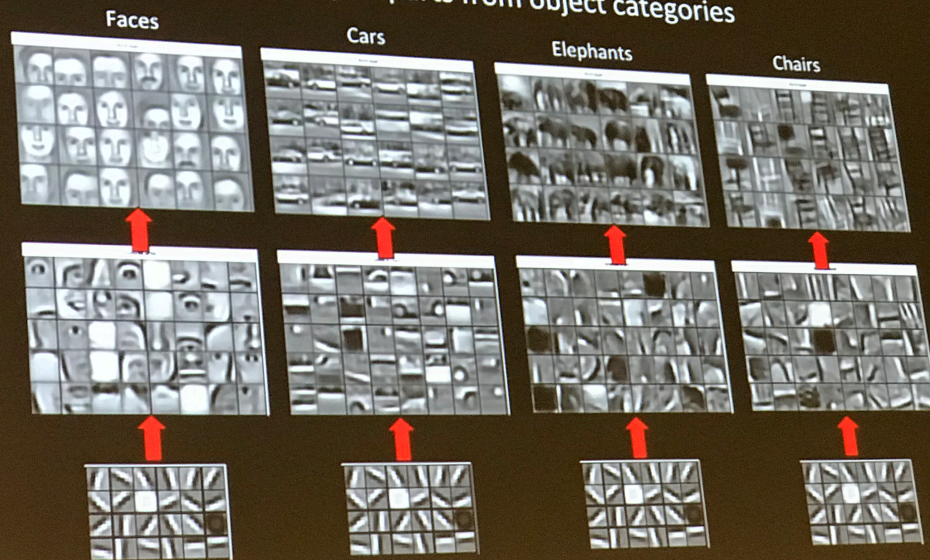
CNN on Face Images 2012



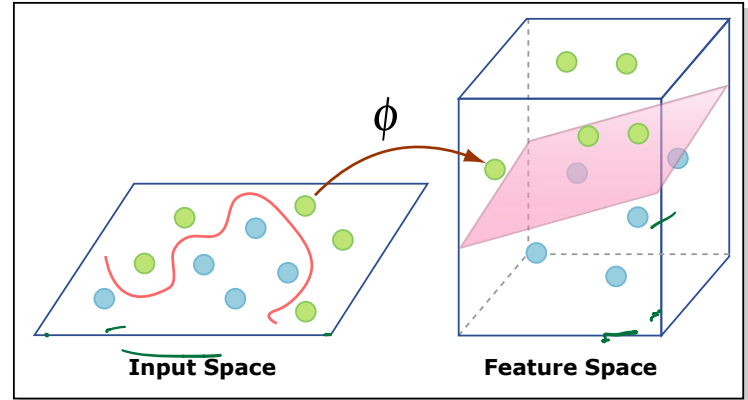
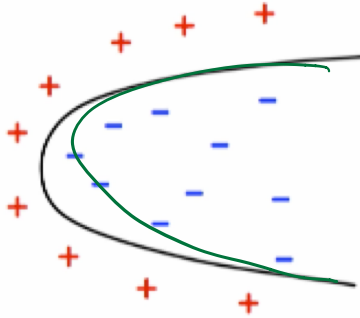
ICASSP 2018 Plenary

ImageNet Large Scale Visual Recognition Challenge, 2012

Examples of learned object parts from object categories



Lecture 9: Kernels

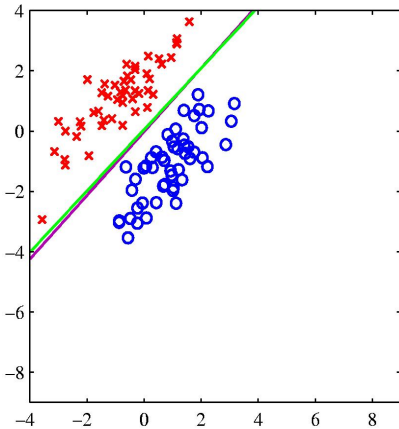


$$w^T \phi(x) \geq 0$$

Basis expansion
Kernel trick

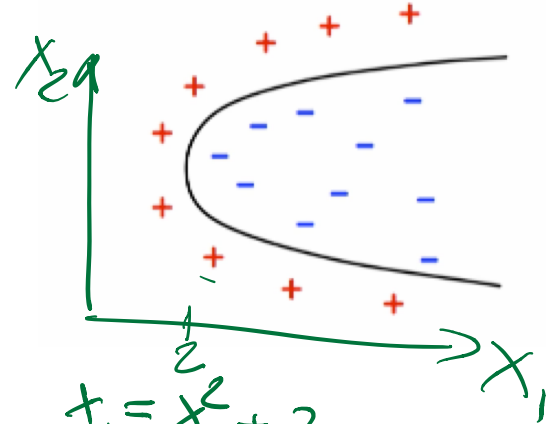
$$M = \infty$$

Basis expansion



$$w^T x = 0$$

$$w^T \phi(x) = 0$$



$$x_1 = x_2^2 + 2$$

$$x_1 - x_2^2 - 2 = 0$$

$$\phi(x) = [1, x_1, x_2, x_1^2, x_2^2, x_1 x_2]$$

$$w = [-2, 1, 0, 0, -1, 0]$$

$$w^T \phi(x) = 0$$

linear in $\phi(x)$

LSQ for classification

Each class \mathcal{C}_k is described by its own linear model so that

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0} \quad (4.13)$$

where $k = 1, \dots, K$. We can conveniently group these together using vector notation so that

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} \quad (4.14)$$

Consider a training set $\{\mathbf{x}_n, \mathbf{t}_n\}, n = 1 \dots N$
Define \mathbf{X} and \mathbf{T}

LSQ solution:

$$\widetilde{\mathbf{W}} = (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} = \widetilde{\mathbf{X}}^\dagger \mathbf{T} \quad (4.16)$$

And prediction

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} = \mathbf{T}^T \left(\widetilde{\mathbf{X}}^\dagger \right)^T \widetilde{\mathbf{x}}. \quad (4.17)$$

Dual representation, Sec 6.2

Primal problem: $\min_{\mathbf{w}} E(\mathbf{w})$

$$\mathbf{X} = \begin{bmatrix} -\mathbf{x}_1^T \\ \vdots \\ -\mathbf{x}_N^T \end{bmatrix} \in \mathbb{R}^{N \times M}$$

$$E = \frac{1}{2} \sum_n^N \{ \mathbf{w}^T \mathbf{x}_n - t_n \}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 = \|\mathbf{X}\mathbf{w} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Solution $\mathbf{w} = \mathbf{X}^+ \mathbf{t} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}_M)^{-1} \mathbf{X}^T \mathbf{t}$ $\mathbf{w} \in \mathbb{R}^M$
 $= \mathbf{X}^T (\mathbf{X} \mathbf{X}^T + \lambda \mathbf{I}_N)^{-1} \mathbf{t} = \mathbf{X}^T (\mathbf{K} + \lambda \mathbf{I}_N)^{-1} \mathbf{t} = \mathbf{X}^T \mathbf{a}$ $\mathbf{a} \in \mathbb{R}^N$

The kernel is $\mathbf{K} = \mathbf{X} \mathbf{X}^T \in \mathbb{R}^{N \times N}$

$$k_{ij} = \phi(\mathbf{x}_i) \phi(\mathbf{x}_j)$$

Dual representation is : $\min_{\mathbf{a}} E(\mathbf{a})$

$$E = \frac{1}{2} \sum_n^N \{ \mathbf{w}^T \mathbf{x}_n - t_n \}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 = \|\mathbf{K}\mathbf{a} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$

Prediction

$$y = \mathbf{w}^T \mathbf{x} = \mathbf{a}^T \mathbf{X} \mathbf{x} = \sum_n^N a_n \mathbf{x}_n^T \mathbf{x} = \sum_n^N a_n k(\mathbf{x}_n, \mathbf{x})$$

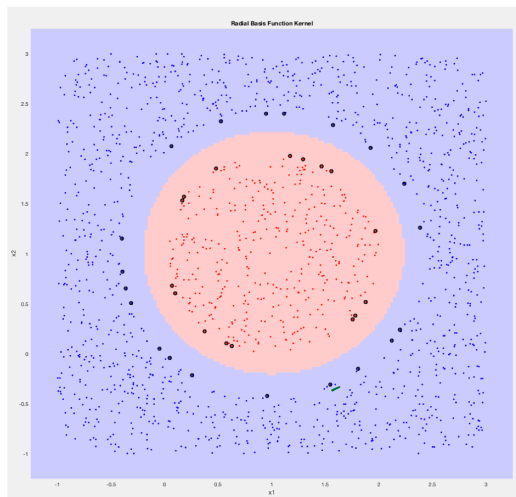
Dual representation, Sec 6.2

Prediction

$$y = \mathbf{w}^T \mathbf{x} = \mathbf{a}^T \mathbf{X} \mathbf{x} = \sum_n^N a_n \mathbf{x}_n^T \mathbf{x} = \sum_n^N a_n k(\mathbf{x}_n, \mathbf{x})$$

- Often **a is sparse** (... Support vector machines)
- We don't need to know **x** or $\varphi(x)$. *Just the Kernel*

$$E(\mathbf{a}) = \|\mathbf{K} \mathbf{a} - \mathbf{t}\|_2^2 + \frac{\lambda}{2} \mathbf{a}^T \mathbf{K} \mathbf{a}$$



Gaussian Kernels

- Gaussian Kernel

$$k(x, x') = \exp \left(-\frac{1}{2} (\mathbf{x} - \mathbf{x}')^T \underline{\Sigma^{-1}} (\mathbf{x} - \mathbf{x}') \right)$$

Diagonal Σ : (this gives ARD)

$$k(x, x') = \exp \left(-\frac{1}{2} \sum_i^N \frac{(x_i - x'_i)^2}{\sigma_i^2} \right)$$

Isotropic σ_i^2 gives an RBF

$$k(x, x') = \exp \left(-\frac{\|\mathbf{x} - \mathbf{x}'\|_2^2}{2\sigma^2} \right) = \phi(x) \phi(x')$$